

I. ПРИКЛАДНЫЕ ВОПРОСЫ ОБУЧЕНИЯ УЗНАВАНИЮ

В. П. Карп, П. Е. Кунин

МЕТОД НАПРАВЛЕННОГО ОБУЧЕНИЯ В ПЕРЕБОРНОЙ СХЕМЕ М. М. БОНГАРДА И ОНКОЛОГИЧЕСКАЯ ДИАГНОСТИКА

Когда М. М. Бонгард выдвинул идею о том, что решение задачи о различении классов объектов следует искать среди возможных логических функций характеристик исходного описания, он сам хорошо понимал, что этот новый подход имеет прямое отношение к медицинской диагностике. В самом деле, стоит только назвать характеристики описания симптомами, как становится ясно, что медицинский термин *синдром* представляет собой не что иное, как бонгардовские конъюнкции характеристик. Понятие синдрома в медицине является очень важным, и в некоторых странах даже классификация заболеваний принята не этиологическая, а синдромная.

М. Бонгард справедливо считал, что узнавание с помощью логических функций характеристик в известной мере моделирует мышление человека; в медицинских задачах это проявляется особенно рельефно.

Конечно, между мышлением человека, решающего задачу узнавания, и мышлением машины, работающей «по Бонгарду», имеется существенное различие. Машина производит полный перебор внутри заданного класса логических функций характеристик описания, отбирая из них самые полезные для узнавания. Человек же, практически не имея возможности осуществлять такой перебор, может только «подбирать» синдромы, кажущиеся ему более существенными для диагностики. Успех такого подбора существенным образом зависит от интуиции, или, иными словами, от накопленного и подсознательно переосмысленного опыта.

Имея в виду только что отмеченное преимущество машины перед человеком, мы начали применять метод Бонгарда для решения одной из самых трудных и самых важных медицинских задач, а именно задачи ранней онкологической диагностики. Слож-

ность этой задачи состоит в том, что симптомы рака, как правило, очень схожи с симптомами тех доброкачественных заболеваний, от которых его нужно отличить.

Метод перебора конъюнкций симптомов (МПК) подкупает не только тем, что как бы имитирует врачебное мышление, но и тем, что он свободен от ряда недостатков, присущих другим методам узнавания.

Вероятностные методы, которые часто используются в медицине (схема Байеса), основываются на предположении о статистической независимости признаков, которое в медицинских задачах мало обосновано. Действительно, наблюдаемые симптомы, как правило, являются проявлением сравнительно небольшого числа нарушений, вызванных заболеванием; кроме того, непосредственный подсчет коэффициентов корреляции между всевозможными парами симптомов почти всегда указывает на то, что они далеко не независимы.

МПК не требует предположений о статистической независимости симптомов и, более того, о характере функций распределения объектов в пространстве рецепторов.

Причина, заставляющая с осторожностью отнестись к применению геометрических методов, связанных с построением разделяющих поверхностей, в задачах медицинской диагностики, заключается в следующем. Если число объектов обучения сравнимо с числом признаков, т. е. с размерностью пространства, то разделяющая гиперплоскость практически может быть найдена всегда (без ошибок в классификации на материале обучения).

В медицинских задачах верифицированный материал обучения обычно весьма ограничен. В то же время число признаков, если не заботиться об их полезности, может быть сделано очень большим. Увеличение же числа признаков при заданном объеме материала обучения всегда дает возможность провести разделяющую гиперплоскость. Ценность такого разделения будет, однако, невелика, если учесть необходимость последующей классификации объектов, не предъявляемых при обучении.

По сравнению с геометрическими методами в МПК следует ожидать меньшего числа вырабатываемых при обучении «предрассудков», так как исследование ведется не в пространстве, размерность которого определяется числом признаков и где плотность точек-объектов обучения обычно мала, а в проекциях этого пространства на подпространства малого числа измерений, где плотность точек материала обучения значительно больше. Кроме того, применение геометрических методов в исходном или в преобразованном пространстве рецепторов предполагает известными значения всех признаков для всех объектов, что в медицинских задачах далеко не всегда осуществимо (например, не все больные проходили одни и те же исследования).

Метод Бонгарда, позволяющий в каждом подпространстве небольшого числа измерений представлять для обучения разное

число объектов, не требует задания всех координат точек, изображающих объекты, т. е. МПК дает возможность использовать в процессе обучения и предъявлять на экзамен объекты с неполным описанием. Однако непосредственное применение разработанного М. Бонгардом и М. Вайнцвайгом алгоритма «Кора-3» к задачам онкологической диагностики оказалось малоэффективным. Причина состоит, в частности, в том, что при различении рака от близких к нему по симптомам заболеваний выявляется очень мало таких конъюнкций, которые встречались бы только в одном из различаемых классов, и все они классифицируют лишь небольшую часть материала обучения, а только такие «чистые» конъюнкции получают отличный от нуля вес в алгоритме «Кора-3».

Конечно, получение отличного от нуля веса только чистыми конъюнкциями не является сколько-нибудь принципиальным ограничением. В общем случае диагностическую оценку или вес можно ввести различными способами. Оказалось, что выбор соотношения, определяющего вес конъюнкций, не сказывается существенным образом на эффективности диагностики. Важно только установить, что в решающий набор допускаются лишь конъюнкции с весом, не меньшим порогового, и для выбранного способа определения веса по материалу обучения найти оптимальное значение порога.

Однако программы типа «Кора» и после указанного расширения класса конъюнкций, допускаемых в решающий набор, продолжали оставаться малоэффективными при диагностике рака. Заметное улучшение было достигнуто после того, как в правило составления решающего набора был введен дополнительный критерий, не связанный с величиной веса конъюнкций.

***k*-ограничение.** Этот дополнительный критерий для построения решающего набора был сформулирован нами следующим образом: решающий набор конъюнкций строится так, чтобы в него входили конъюнкции с возможно большим весом (во всяком случае, превышающим пороговый) и чтобы при этом ни один из исходных симптомов не входил более чем в *k* конъюнкций решающего набора. Последнее условие мы назвали *k*-ограничением. Число *k* является дискретным параметром алгоритма. Оптимальными значениями параметра *k* оказались небольшие числа: $k = 1-3$.

Причина диагностической эффективности *k*-ограничения, по видимому, заключается в том, что оно существенно лимитирует возможность попадания в решающий набор конъюнкций, между которыми на материале обучения наблюдаются значительные корреляции.

Идея *k*-ограничения перекликается с высказыванием М. Бонгарда о том, что конъюнкции решающего правила должны быть по мере возможности независимы. Действительно, *k*-ограничение косвенным образом не допускает в решающий набор сильно коррелированные между собой конъюнкции. Конечно, такой способ учета корреляций является очень грубым. По существу здесь

конъюнкции, имеющие общий симптом, предполагаются сильно зависимыми, в противном случае они считаются независимыми или слабо зависимыми. Нетрудно привести примеры, показывающие, что, вообще говоря, это не так. Тем не менее даже столь грубый способ защиты решающего набора от проникновения зависимых конъюнкций увеличивает эффективность диагностики.

Представляется, конечно, очень заманчивым учесть корреляции между конъюнкциями не косвенным, а прямым путем. Однако в рамках традиционного подхода к задаче узнавания это оказывается невозможным.

Направленное обучение. При обычном подходе к задаче узнавания найденное в процессе обучения решающее правило фиксируется «раз и навсегда», т. е. в дальнейшем применяется ко всем возможным контрольным объектам. Поэтому естественно говорить, что в таких случаях имеет место универсальное обучение. Почти все применяемые алгоритмы узнавания предполагают универсальность процесса обучения.

Нами был предложен новый подход к решению задачи узнавания. Предлагается процесс обучения проводить направленно — отдельно для каждого контрольного объекта — таким образом, чтобы изучались зависимости только между симптомами, которые имеются у представленного на экзамен объекта. В такой постановке задачи предпосылкой процесса обучения является задание значений признаков (симптомов) не только всех объектов материала обучения (для каждого из которых известно, к какому из различаемых классов он относится), но также и представленного на экзамен объекта, принадлежность которого к одному из классов (A или B) предстоит определить.

Такое обучение мы назвали направленным обучением, имея в виду, что выработка решающего правила происходит «направленно» по отношению к каждому объекту, представленному на экзамен. При направленном обучении с использованием, например, МПК, все симптомы, вошедшие в конъюнкции решающего набора, заведомо имеются у экзаменуемого объекта.

Формально направленное обучение состоит в следующем. Пусть объекты, правило классификации которых требуется найти, характеризуются n признаками, причем s -й признак имеет l_s допустимых реализаций (симптомов). Пусть, далее, $\sum_{s=1}^n l_s = N$. Тогда

каждый объект X есть точка в N -мерном пространстве R_N с координатами X^i ($i = 1, 2, \dots, N$), принимающими значения 0 или 1, причем число координат, равных единице, есть n , если для данного объекта известны значения всех признаков. Допустим и случай, когда известны значения не всех признаков, тогда число ненулевых координат будет меньше n .

Пусть T — данный объект и i_1, i_2, \dots, i_m ($m \leq n$) — номера тех координат, для которых $t^{i_s} = 1$. Через $R_m(T)$ обозначим

m -мерное подпространство пространства R_N , натянутое на оси с номерами i_s .

Обозначим через P_T оператор проектирования пространства R_N в подпространство $R_m(T)$. Заметим, что координаты проекции $P_T(X)$ объекта $X \in R_N$ в пространство $R_m(T)$ получаются просто отбрасыванием всех координат X^i с $i \neq i_s$ ($s = 1, 2, \dots, m$).

Легко видеть, что любое решающее правило (хотя бы и полученное при рассмотрении объектов в пространстве R_N) может быть применено к объекту T в пространстве $R_m(T)$, так как единичное значение координаты, соответствующей симптому некоторого признака, предопределяет нулевые значения остальных координат, порожденных тем же признаком.

Пусть a_1, a_2, \dots, a_p и b_1, b_2, \dots, b_q — объекты наборов обучения, принадлежащие соответственно классам A и B . Тогда $P_T(a_k)$ ($k = 1, 2, \dots, p$) и $P_T(b_k)$ ($k = 1, 2, \dots, q$) — проекции этих объектов в $R_m(T)$. Переход к направленному обучению состоит в построении решающего правила для отнесения объекта T к одному из двух классов A и B с использованием объектов $P_T(a_k)$ и $P_T(b_k)$ в качестве набора обучения.

Любой алгоритм, применяемый к $\{P_T(a_k)\}$ и $\{P_T(b_k)\}$ в целях классификации объекта T , называется алгоритмом направленного обучения. В принципе любой алгоритм, действующий в произвольном пространстве дихотомических признаков, может быть использован в качестве алгоритма направленного обучения, так как пространство $R_m(T)$ может рассматриваться в качестве исходного. Фактически дело обстоит сложнее. Так, например, алгоритм распознавания с помощью проведения разделяющей гиперповерхности может отказать при переходе от R_N к $R_m(T)$, так как проекции объектов обучения классов A и B могут не разделяться, скажем, гиперплоскостью в $R_m(T)$, в то время как сами объекты разделялись гиперплоскостью в R_N . Правда, в этом случае можно ставить задачу по-другому: например, можно искать гиперплоскость в $R_m(T)$, наилучшим образом отделяющую проекции объектов класса A от проекций объектов класса B , но это уже по существу меняет постановку задачи.

При направленном обучении претерпевает изменение основной принцип поиска решающего правила. Действительно, при универсальном обучении построение решающего правила основывается на стремлении как можно лучше разделить материал обучения, т. е. получить максимальное возможное число правильных классификаций при минимальном числе ошибок (при этом можно учитывать и то, что два возможных типа ошибок имеют разную цену).

При направленном обучении такой подход не всегда осуществим хотя бы потому, что ряд объектов обучения обоих классов в результате проектирования может попасть в одни и те же точки подпространства $R_m(T)$, и их разделение становится принципиально неосуществимым. Тем не менее ни один из симптомов, которыми различаются упомянутые объекты обучающегося набора,

не относится к симптомам объекта T , и поэтому такие «чужие» для него симптомы вряд ли несут полезную информацию для его диагностики. При направленном обучении помимо алгоритмов, индуцированных алгоритмами универсального обучения, возникает возможность создания других, специфических алгоритмов, которые не переносятся (по крайней мере, непосредственно) на универсальное обучение.

Так, в рамках направленного обучения оказалось возможным создать алгоритм, приводящий к решающему набору, который содержит только такие конъюнкции (с возможно большими весами), что коэффициенты корреляции между любыми двумя из них не превышают некоторого заданного порогового значения. Как будет видно из самого строения алгоритма, его эффективная реализация возможна только при направленном обучении. Этот алгоритм предусматривает построение решающего набора в два этапа. На первом этапе для экзаменуемого объекта создаются списки конъюнкций, являющихся «кандидатами» в решающий набор, причем в них попадают все конъюнкции, имеющие вес, больший порогового. На втором этапе из конъюнкций-кандидатов строится решающий набор. В каждом из двух кандидатских списков (соответствующих возможным диагнозам A и B) конъюнкции упорядочиваются в соответствии с их весами. Первая конъюнкция (с наибольшим весом) заносится в решающий набор без проверки. Затем производится подсчет коэффициентов корреляции между этой занесенной в решающий набор конъюнкцией и всеми остальными конъюнкциями кандидатского списка. Те конъюнкции, коэффициент корреляции которых с первой превышает заданный порог, вычеркиваются из списка кандидатов, а наиболее весомая из оставшихся конъюнкций переводится в решающий набор. По отношению к вновь включенной в решающий набор конъюнкции производится такая же процедура, и процесс этот повторяется до тех пор, пока в списке кандидатов не останется ни одной конъюнкции. Таким путем составляются обе части решающего набора.

Применения в онкологической диагностике. С помощью метода перебора конъюнкций, используя описанный алгоритм направленного обучения, мы исследовали ряд задач альтернативной онкологической диагностики. Результаты решения четырех из них приводятся ниже. Эти четыре задачи таковы:

а) альтернативная рентгенодиагностика между периферическим раком (ПР) и доброкачественными образованиями (ДО) в легких (задача ПР — ДО, материал обучения — 261, контрольная группа — 100);

б) альтернативная рентгенодиагностика между центральным раком (ЦР) и неспецифическими хроническими воспалительными процессами (ХВП) в легком (задача ЦР — ХВП, материал обучения — 287, контрольная группа — 100);

в) альтернативная диагностика между раком и доброкачественными заболеваниями молочной железы по клиническим и ци-

тологическим данным (задача РМЖ — ДЗ, материал обучения — 131, контрольная группа — 100);

г) альтернативная диагностика между злокачественной и доброкачественной язвой желудка по клиническим и рентгенологическим данным (задача ЗЯ — ДЯ, материал обучения — 118, контрольная группа — 60).

Во всех этих задачах алгоритмы направленного обучения оказались эффективнее алгоритмов универсального обучения и значительно эффективнее схемы Байеса, основанной на предположении о статистической независимости признаков исходного описания.

Результаты диагностики с помощью различных машинных алгоритмов на одной и той же группе контрольных больных в сравнении с клинической диагностикой тех же больных представлены в таблице.

Метод диагностики	Относительное количество диагнозов, %			Метод диагностики	Относительное количество диагнозов, %		
	правильных	неопределенных (отказов от диагноза)	ошибочных		правильных	неопределенных (отказов от диагноза)	ошибочных
ПР — ДО				РМЖ — ДЗ			
Диагнозы рентгенологов	66	24	10	Диагнозы цитологов	78	17	5
Схема Байеса	84	8	8	Схема Байеса	90	6	4
Универсальное обучение	84	11	5	Универсальное обучение	94	4	2
Направленное обучение	92	4	4	Направленное обучение	95	3	2
ПР — ХВП				ЗЯ — ДЯ			
Диагнозы рентгенологов	62	23	15	Клинические диагнозы	21	54	25
Схема Байеса	72	18	10	Схема Байеса	67	12	21
Универсальное обучение	80	15	5	Универсальное обучение	77	13	10
Направленное обучение	85	9	6	Направленное обучение	83	7	10